# INFORMATION RETRIEVAL FOR NONSTATIONARY DATA RECORDS

FINAL REPORT

Prepared for

## NATIONAL AERONAUTICS AND SPACE ADMINISTRATION
## GEORGE C. MARSHALL SPACE FLIGHT CENTER
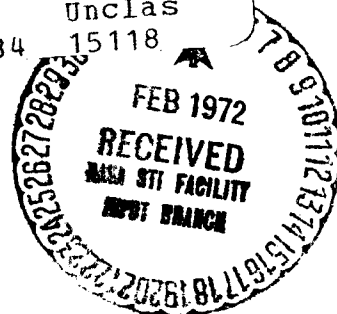### Aero-Astrodynamics Laboratory

Under Contract NAS8-26757

**NORTHROP SERVICES. INC.**

P. O. BOX 1484

HUNTSVILLE, ALABAMA  35807

TELEPHONE (205)837 0580

FINAL REPORT

# INFORMATION RETRIEVAL FOR NONSTATIONARY DATA RECORDS

December 1971

By

M. Y. Su

PREPARED FOR:

*NATIONAL AERONAUTICS AND SPACE ADMINISTRATION*
*GEORGE C. MARSHALL SPACE FLIGHT CENTER*
*AERO-ASTRODYNAMICS LABORATORY*

*Under Contract NAS8-26757*

REVIEWED AND APPROVED BY:

*A. L. Grady*

A. L. Grady, Manager
Advanced Engineering Analysis

**NORTHROP SERVICES, INC.**
**HUNTSVILLE. ALABAMA**

# FOREWORD

This study was undertaken by Northrop Services, Inc., for NASA/MSFC, Huntsville, Alabama, under Contract No. NAS8-26757. The program was under the direction of MSFC Aero-Astrodynamic Laboratory, Flight Data Statistics Office, with Mr. R. R. Jayroe, Jr. as the project monitor.

# ACKNOWLEDGEMENTS

The author wishes to acknowledge Mr. R. R. Jayroe, Jr. for his contribution to Section II on the Envelope Algorithm for Splitting Nonstationary Time Series. Mr. Jayroe proposed the algorithm and also made the application to the sunspot data.

## SUMMARY

Most random time series recorded in the natural (uncontrolled) environment such as in the atmosphere are affected by the variations of the environment. These time series are statistically nonstationary. This is particularly true for remote observations of atmosphere trace constituents, since the radiometers may see several layers of constituents over a large range of altitudes.

While a very extensive theory exists on the analysis of stationary random processes, there is no generally applicable theory and methodology for analysis of nonstationary time series. In Section I, a review and critical discussion will be made on the existing methods for analysis of nonstationary time series. In Section II, a new algorithm for splitting nonstationary time series, which was originally proposed by R. R. Jayroe, Jr., will be presented and discussed. This algorithm has been applied to analysis of the sunspot data. Finally, some conclusions will be drawn in Section III.

# TABLE OF CONTENTS

# LIST OF ILLUSTRATIONS

# LIST OF ILLUSTRATIONS (Concluded)

# Section I

# METHODS FOR ANALYZING NONSTATIONARY DATA

Generally speaking, there are three types of approaches to the analysis of nonstationary random time series. The first follows the classical theory of stationary time series; statistical quantities such as mean, correlation and spectrum are defined much the same ways as for stationary time series except that an additional time or frequency parameter is introduced. The formalism of the approach may be similar to stationary time series analysis, yet the physical interpretation of such double-parameters quantities is a critical issue and not without ambiguity (ref. 1).

The second approach follows the model construction method by assuming that the nonstationary time series consists of a slow-varying trend, a periodic component and stationary random component. The basic strategy of this approach is to split the nonstationary time series into its components, and then analyze each component separately (ref. 2).

It is interesting to note that both above-mentioned references appear in the same year, 1966, and that both are supposed to present general methods for analysis of nonstationary time series. However, their treatments are exclusive of each other; reference 1 discusses only the classical formal approach, while reference 2 only the trend-elimination approach. This might indicate a degree of subjectivity in dealing with nonstationary processes.

The third approach is a further refinement of the first approach, in that an optimum filtering operation is derived and performed on the double-parameters correlation or spectrum function obtained from the ensemble average (ref. 3). Two test functions of error measures needed to be assumed, which will be minimized to obtain its optimum smoothing feature. The whole formalism seems to be too complicated to be of practical application. For the latter reason, no more discussion will be made on this approach.

It may be fair to say that the classical ensemble approach is easier to process, but harder to interpret the obtained results physically. On the other

hand, the trend-elimination approach is harder to process, but easier to interpret the obtained results. In particular, when there exists only one or two time series, the later approach has to be adopted. Clearly, either approach is far from satisfactory for general analysis of nonstationary time series. Even for dealing particular time series, meaningful results may still require some ingenuity and some a priori knowledge from the investigator.

A piecewise detrending method and piecewise weighting method have also been proposed and applied to the crossed beam data with partial success (refs. 4 through 6). Since these later two methods have been discussed in detail in the cited references, no more discussions will be made in this report.

In this section, the first two approaches and their merits and disadvantages will be reviewed and discussed. Subsection 1.1 presents the formal classical approach and subsection 1.2 presents the trend-elimination approach.

## 1.1 THE ENSEMBLE APPROACH

Let $X_i(t)$ and $Y_i(t)$; $0 \leq t \leq T$; $i = 1, 2..., N$ denotes N pairs of nonstationary random time series. The mean, variance, covarince, and first-order probability density functions are then defined by the ensemble averaging as follows.

Mean:
$$\hat{\mu}_x(t) = \frac{1}{N} \sum_{i=1}^{N} X_i(t) \tag{1}$$

Variance:
$$\hat{\sigma}_x^2(t) = \frac{1}{N} \sum_{i=1}^{N} [X_i(t) - \hat{\mu}_x(t)]^2. \tag{2}$$

Covariance:
$$\hat{R}_{xy}(t_1,t_2) = \frac{1}{N} \sum_{i=1}^{N} [X_i(t_1) - \hat{\mu}_x(t_1)] [Y_i(t_2) - \hat{\mu}_y(t_2)] \tag{3}$$

First-order probability density:

$$\hat{p}(X(t)) = \frac{1}{\Delta X} \quad [X - \Delta X \leq X(t) \leq X + \Delta X] \tag{4}$$

With the transformation

$$\tau = t_2 - t_1$$

$$t = \frac{1}{2} (t_1 + t_2) \quad , \tag{5}$$

the covariance can be rewritten as

$$\hat{R}_{xy} (t_1, t_2) = \hat{R}_{xy} (t - \frac{\tau}{2}, t + \frac{\tau}{2}) \tag{6}$$

As for spectral densities, there are four different definitions as follows;

(a)  Double frequency spectra

(b)  Time-varying power spectra

(c)  Time-averaging power spectra

(d)  Instantaneous (frequency – time) spectra.

Definitions (a) and (d) are considered to be significant theoretical ways to analyze nonstationary spectra. Definitions (b) and (c) are experimentally measurable by direct filtering procedures. Each of these spectra will be defined next.

### 1.1.1 Double Frequency Spectrum

Double Fourier transform of the two-times covariance $\hat{R}_{xy} (t_1, t_2)$ yields the double frequency spectrum

$$\hat{S}_{xy} (f_1, f_2) = \int\int_{-\infty}^{\infty} \hat{R}_{xy} (t_1, t_2) \exp[-j2\pi(f_1 t_1 - f_2 t_2)] \, dt_1 \, dt_2$$

$$= E [X_i^* (f_1) Y_i (f_2)] \tag{7}$$

where

$$X(f_1) = \int_{-\infty}^{\infty} X_i(t) \exp[-j2\pi f_1 t] \, dt \tag{8}$$

and the asterisk denotes the complex conjugate. In general, the double frequency spectum is a complex function. With the following change of variables.

$$f = \frac{1}{2} (f_1 + f_2)$$

$$g = f_2 - f_1 \tag{9}$$

we have

$$\hat{S}_{xy}(f,g) = \hat{S}_{xy}^{(1)}(f_1,f_2)$$

$$= \int\int_{-\infty}^{\infty} \hat{R}_{xy}(t,\tau) \exp[-j2\pi(f\tau + gt)] \, dt \, d\tau \tag{10}$$

### 1.1.2 Time-Varying Power Spectrum

Each of the random time series, $X_i(t)$, $i = 1,2...,N$, is filtered through a narrow-bandpass filter with central frequency f and bandwidth B, and then the instantaneous outputs are squared. Averaging the N outputs and dividing by B gives the time-varying power spectrum. Mathematically, these operations can be expressed as follows:

$$X_i(t, f, B) = \int_{-\infty}^{\infty} h(\tau) \, X_i(t - \tau) \, dt \tag{11}$$

where $h(\tau)$ is the transfer function of the filter whose frequency response function is given by

$$H_f(f') = 1 \qquad\qquad f' - \frac{B}{2} \leq f \leq f' + \frac{B}{2}$$

$$= 0 \qquad\qquad \text{otherwise.} \tag{12}$$

Next,

$$X_i^2(t, f, B) = \int\int_{-\infty}^{\infty} h(\tau_1) \, h(\tau_2) \, X_i \, (t - \tau_1) \, X_i \, (t - \tau_2) \, d\tau_1 \, d\tau_2 \tag{13}$$

Finally, the time-varying power spectrum is

$$\hat{G}_x(t, f, B) = \frac{1}{BN} \sum_{i=1}^{N} X_i^2 \, (t, f, B) \tag{14}$$

which is positive for all values of t and f. If the number of time series, N, approaches infinite, then the "expected" time-varying power spectrum

$$G_x(t, f) = \frac{1}{B} \lim_{N \to \infty} \frac{1}{N} \sum_{i=1}^{N} X_i^2 \, (t, f, B)$$

$$= \frac{1}{B} \, E \, [X_i^2 \, (t, f, B)] \tag{15}$$

is obtained. This is related to the double frequency spectrum and double-time covariance;

$$G_x(t,f) = \frac{1}{B} \int\int_{-\infty}^{\infty} H^*(f_1) \, H(f_2) \, S_x(f_1, f_2) \, \cdot$$

$$\exp[-j2\pi(f_1 - f_2)t] \, df_1 \, df_2$$

$$= \frac{1}{B} \int\int_{-\infty}^{\infty} h(\tau_1) \, h(\tau_2) \, R_x \, (t - \tau_1, t - \tau_2) \, d\tau_1 \, d\tau_2 \tag{16}$$

### 1.1.3 Time—Averaging Power Spectra

This spectrum is defined as the time average over T of the time-varying power spectrum

$$\bar{G}_x \, (f, t, T) = \frac{1}{T} \int_{t - T/2}^{t + T/2} G_x(t', f) \, dt' \tag{17}$$

which is non-negative.  It can be rewritten as

$$\bar{G}_x(f, t, T) = \frac{1}{B} \int\int_{-\infty}^{\infty} h(\tau_1) \, h(\tau_2) \, [\frac{1}{T} \int_0^T R_x(t - \tau_1, \, t - \tau_2) dt] \, d\tau_1 \, d\tau_2$$

$$= \frac{1}{B} \int\int_{-\infty}^{\infty} h(\tau_1) \, h(\tau_2) \, \bar{R}_x(\tau_1 - \tau_2; t) \, d\tau_1 \, d\tau_2 \qquad (18)$$

where

$$\bar{R}_x(\tau_1 - \tau_2; \, t) = \frac{1}{T} \int_0^T R_x(t - \tau_1, \, t - \tau_2) \, dt.$$

The averaging time T can only be determined by trial and error such that T is long enough to smooth out the instantaneous fluctuation, but short enough not to introduce significant bias error which reflects the smoothing of nonstationary trends in the data.

### 1.1.4 Instantaneous (Frequency-Time) Spectra

This spectrum is obtained by a single Fourier transform of the double-time covariance with respect to only the time lag $\tau$;

$$S'_{xy}(f, t) = \int_{-\infty}^{\infty} \hat{R}_{xy}(t, \tau) \, \exp(-j2\pi f\tau) \, d\tau \qquad (19)$$

This is called instantaneous, since the spectrum is associated with any instant of time t.  Actually, it is a function of frequency and time simultaneously.

From the above definitions and discussions, it can be seen that this ensemble approach for analyzing nonstationary time series is not applicable to remote detection of the atmospheric air constituents because there is only one pair of time series with sufficient length of time to be dealt with.

## 1.2 TREND-ELIMINATION APPROACH

By the trend elimination approach, a random nonstationary time series is assumed to consist of three independent components;

- Trends
- Periodic components
- Stationary random time series.

The strategy of the approach is to decompose the given time series into each individual component and then analyze each separately. In fact, many practical problems may be quite realistically represented by such a model. This fact may be the best justification for this approach.

The concept of a trend does not have a clear-cut definition. Generally, one thinks of it as a smooth change of the phenomenum under consideration over a long period of time, but "long" in this connection is a relative term. For practical purposes, the length of trends should be at least one order of magnitude larger than the period of the regular oscillation and random fluctuations.

The concept of oscillations with a constant period is much easier to grasp and more definite. Seasonal changes with a period of one year and daily change with 24 hours period are the two most notable and common examples in natural atmospheric environments. Because the period of these regular oscillations is known, it will be easier to isolate this component from the time series.

However long a time series may be, one can never be certain and often not even reasonably sure that a trend in it is not part of a slow oscillation, except of course when the time series has terminated for a practical reason or limited length of measurement. In speaking of a trend, one should bear in mind the length of the time series to which our statement refers. Perhaps it would be more accurate to speak of slow or quick movements rather than of trends and oscillations, but even so the distinction between these two would still remain a matter of subjective judgement to some extent.

### 1.2.1 Determination of Trend

Since the trend, by the above general description, is smooth over a fairly long interval of time, it is reasonable to assume that, at least, within a finite interval around any time t, the trend may be approximated by a poly-nominal in time.  Thus, given a time series X(t), an $m^{th}$ degree polynominal may be sought

$$g(t) = \sum_{i=0}^{m} a_i t^i \tag{20}$$

where the coefficient $a_i$ will be determined by using the time series.

If a polynominal is fitted to the whole time series by the least square method, it gives the curvilinear regression line of X(t).  However, to obtain a satisfactory trend curve, one would need to have a very high-degree polynominal. Obviously, this is impractical for most time series with fairly complicated trend curves.  An alternative is to fit a polynominal only to any small part of the time series.  The simplest method which also forms the basis of the majority of the trend elimination method is to use a polynominal of degree m to consecutive (2n + 1) digital samples, which are taken at an interval of $\Delta t$ and also m<2n.

The coefficients of the polynominal will be determined by the least square method, i.e., by solving the following system of equations

$$\frac{\partial}{\partial a_i}\left\{ \sum_{j=-n}^{n} [X_{k+j} - \sum_{i=0}^{m} a_i (j\Delta t)^i]^2 \right\} = 0 \tag{21}$$

i = 0, 1, 2,...,m and any integer k.
This yields a system of (m + 1) linear equations with (m + 1) unknowns, $a_i$. Here $X_k$ denotes the sample X(k$\Delta$t).  Since the middle term is the only one needing evaluated, which is equal to $a_o$ at t = 0, the trend will be of the form

$$g(t) = g(k\Delta t) = a_o = \sum_{j=-n}^{n} W_j (m,n) X_{k+j} \tag{22}$$

where $W_j (m,n)$ is the function of m and n only.  As can be seen equation (22) is nothing but a weighted moving average of the given time series and the

weights are independent of the time series. Further, the weights are symmetric in n, i.e.

$$W_j(m,n) = W_{-j}(m,n) \qquad\qquad j = 0, \pm1, \pm2\ldots\pm n \qquad (23)$$

It can be further proved that the weights for a polynominal of 2m degree and (2m+1) degree are the same, i.e.,

$$W_j(2m,n) = W_j(2m+1, n) \qquad\qquad j = 0, \pm1, \pm2\ldots,\pm n \qquad (24)$$

Tables 1 and 2 (pages 368 and 369, ref. 2) list the weights for polynominals of degree 2 to 5 with number of samples used being from 5 to 21.

One drawback of the moving average method is its failure to provide the trend values for the first n and the last n samples of the original time series. However, it is not a great loss, if the time series is long enough.

### 1.2.2 Effect of Trend—Elimination By The Moving Average On Other Components

The effect of the moving average on oscillating and random components in the process of removing the trend will now be discussed.

Let the time series $X(t)$ consisting of three components; a trend $X_1(t)$, an oscillatory term with a regular period $T_o$. $X_2(t)$, and a stationary random component $X_3(t)$, i.e.,

$$X(t) = X_1(t) + X_2(t) + X_3(t) \qquad (25)$$

Let the operation of the moving average for trend removed be denoted by $T[\ ]$, then

$$T[X(t)] = T[X_1(t)] + T[X_2(t)] + T[X_3(t)] \qquad (26)$$

It will now be assumed that the method of determining the trend is completely correct in the sense that

$$T[X_1(t)] = X_1(t) \qquad (27)$$

Table 1. WEIGHTS FOR SECOND AND THIRD DEGREE POLYNOMINALS APPROXIMATION OF TRENDS

| j | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---|---|---|---|---|---|---|---|---|---|----|----|
| 35 $W_j(3,5)$ | 17 | 12 | -3 | | | | | | | | |
| 21 $W_j(3,7)$ | 7 | 6 | 3 | -2 | | | | | | | |
| 231 $W_j(3,9)$ | 59 | 54 | 39 | 14 | -21 | | | | | | |
| 429 $W_j(3,11)$ | 89 | 84 | 69 | 44 | 9 | -36 | | | | | |
| 143 $W_j(3,13)$ | 25 | 24 | 21 | 16 | 9 | 0 | -11 | | | | |
| 1105 $W_j(3,15)$ | 167 | 162 | 147 | 122 | 87 | 42 | -13 | -78 | | | |
| 323 $W_j(3,17)$ | 43 | 42 | 39 | 34 | 27 | 18 | 7 | -6 | -21 | | |
| 2261 $W_j(3,19)$ | 269 | 264 | 249 | 224 | 189 | 144 | 89 | 24 | -51 | -136 | |
| 3059 $W_j(3,21)$ | 329 | 324 | 309 | 284 | 249 | 204 | 149 | 84 | 9 | -76 | -171 |

Table 2. WEIGHTS FOR FOURTH AND FIFTH DEGREE POLYNOMINAL APPROXIMATION OF TRENDS

| j | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---|---|---|---|---|---|---|---|---|---|----|----|
| 231 $W_j$(5,7) | 131 | 75 | -30 | 5 | | | | | | | |
| 429 $W_j$(5,9) | 179 | 135 | 30 | -55 | 15 | | | | | | |
| 429 $W_j$(5,11) | 143 | 120 | 60 | -10 | -45 | 18 | | | | | |
| 2431 $W_j$(5,13) | 677 | 600 | 390 | 110 | -135 | -198 | 110 | | | | |
| 46189 $W_j$(5,15) | 11063 | 10125 | 7500 | 3755 | -165 | -2937 | -2860 | 2145 | | | |
| 4199 $W_j$(5,17) | 883 | 825 | 660 | 415 | 135 | -117 | -260 | -195 | 195 | | |
| 7429 $W_j$(5,19) | 1393 | 1320 | 1110 | 790 | 405 | 18 | -290 | -420 | -255 | 340 | |
| 260015 $W_j$(5,21) | 44003 | 42120 | 36660 | 28190 | 17655 | 6378 | -3940 | -11220 | -13005 | -6460 | 11628 |

Subtracting the above two equations to remove the trend,

$$Y(t) \equiv X(t) - T[X(t)]$$

$$= X_2(t) + X_3(t) - T[X_2(t)] - T[X_3(t)] \tag{28}$$

Clearly, the last two terms denotes the effect of the moving average on the oscillating and random components. Consider first $T[X_2(t)]$. To analyze this term in general is quite difficult, so a simple case will be considered where $X_2(t)$ is a sine wave and the weights of the moving average are constant. Hence,

$$T[X_2(t)] = T[X_2(i\Delta t)]$$

$$= \frac{1}{2n+1} \sum_{j=-n}^{n} \sin[f_o + \Delta f(i-j)]$$

$$= \left\{ \frac{1}{2n+1} \quad \frac{\sin \frac{1}{2}(2n+1)\,\Delta f}{\sin \frac{1}{2}\,\Delta f} \right\} \sin[f_o + (n+1)\,\Delta f] \tag{29}$$

which is a sine wave of the same frequency and phase as the orignal, but with the amplitude reduced by the factor shown in the paraenthese { }. Thus, the term $T[X_2]$ will be small if n is large, or if 1/2 (2n+1) $\Delta f$ is a multiple of $\pi$, i.e., if the length of the moving average equals a multiple of the period of the oscillation. On the other hand, if both n and $\Delta f$ are small, then the factor of reduction of amplitude will be approximately equal to unity. This implies then that $T[X_2(t)] \simeq X_2(t)$. Hence, the moving average has eliminated the oscillating component as well as the trend. This result is expected, since a slow ocillation is treated as a trend by the moving average and eliminated accordingly. Generally, the moving average will emphasize the shorter oscillations at the expense of the longer ones. It is not easy to exhibit the precise effect of the moving average when the weights are unequal and the $X_2(t)$ is not harmonic, but evidently the same kind of situation is apt to arise.

Next, consider the effect of the simple average on the random component $X_3(t)$, i.e.,

$$T[X_3(t)] = T[X_3(i\Delta t)]$$

$$= \frac{1}{2n+1} \sum_{j=-n}^{n} X_3[(i + j)\, \Delta t] \qquad (30)$$

It is noted that consecutive terms of $X_3(i\Delta t)$ are less correlated than the consecutive terms of $T[X_3(i\Delta t)]$, since $T[X_3(i_1\Delta t)]$ and $T[X_3(i_2\Delta t)]$ have $(2n+1) - (i_2 - i_1)$ terms of $X_3(t)$ in common, if $(i_2 - i_1) < 2n+1$. Thus, the effect of the moving average is to introduce some artificial correlation or coherence to the random component.

In summary, the effect of taking a moving average of random time series will then be to generate an oscillatory time series, provided the weights are such as to give a positive correlation between successive members of the generated series, a condition which is always realized in moving average for trend-fitting.

Pugachev (ref. 7) gave a more rigorous derivation of the condition for the applicability of the straight moving average to the effect that the random function should be approximately linear over the period $T_o$ of the moving average and the mean value of the covariance over a square with sides of $T_o$ around the mid-point $(t,t)$ should be small. It would be expected that similar conditions should be also imposed on the weighted moving average.

One method for reducing effects of the moving average on components other than the trend may be described as follows: Apply the trend-elimination operator T repeatedly, say m times, to the output of each preceding operation, one then obtains

$$T^m[X] = T^m[X_1] + T^m[X_2] + T^m[X_3] \qquad (31)$$

Assume that the method of determining the trend is completely correct in the sense that

$$T[X_1] = X_1$$

then

$$T^m[X] = X_1 + T^m[X_2] + T^m[X_3] \qquad (32)$$

Now, it has been shown that the amplitude of $T[X_2]$ is reduced by a factor $R \leq 1$ for each operation, and that the expected standard derivation of $T[X_3]$ is reduced by a factor of $(\sqrt{2m+1})^{-1}$. Hence, after a sufficient number of moving averages, the last two terms may be made as small as desired and

$$T^n[X] = X_1 \text{ for m large}$$

and

$$Y(t) = X(t) - T^m[X] = X_2(t) + X_3(t). \qquad (33)$$

Next, there is a need to separate the remaining two components $X_2(t)$ and $X_3(t)$ from $Y(t)$. If the period of $X_2(t)$ is known as in many practical cases, then this separation is fairly easy. Let the known period be denoted by $T_o$, and break the time series $Y(t)$ into shorter pieces with length $T_o$. That is,

$$\{Y_i(t)\} \qquad\qquad (i - 1)T_o \leq t \leq i\,T_o,$$

with $i = 1, 2, \ldots, M \leq T/T_o$

Take the ensemble average of $\{Y_i\}$

$$\overline{Y_{T_o}}(t) = \frac{1}{M} \sum_{i=1}^{M} Y_i(t)$$

$$= \frac{1}{M} \sum_{i=1}^{M} [X_{2,i}(t) + X_{3,i}(t)]$$

$$= X_2(t) + \frac{1}{M} \sum_{i=1}^{M} X_{3,i}(t) \qquad (34)$$

For a large M, the last term becomes small. Thus,

$$\overline{Y_{T_o}}(t) = X_2(t)$$

and $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad (35)$

$$X_3(t) = Y(t) - \overline{Y_{T_o}}(t)$$

# Section II

# ENVELOPE ALGORITHM FOR SPLITTING TIME SERIES

The envelope algorithm for splitting nonstationary time series will be described first and then discuss its characteristics. This algorithm was originally proposed by R. R. Jayroe, Jr. Next, results of application of the algorithm to the sunspot data will be presented.

## 2.1 THE ENVELOPE ALGORITHM

The use of the concept of an envelope function of a time series in communication engineering is not new and has been studied, for example, by R. Deutsch (ref. 8). However, the use of an envelope as proposed by R. Jayroe is different from that treated in the given reference. The reason for using the envelope algorithm stems from the following argument. If one plots and examines the time series to be worked with, it is possible to draw a line through the data about which the data fluctuates. If the data values of this line are subtracted from the original data, then the data will oscillate approximately about a zero value and no long term trends will be present. The envelope algorithm attempts to approximate such a line.

The envelope algorithm may be described as follows: let the given time series under consideration be $x(t)$, $t \geq 0$. Connecting any two consecutive peaks (or local maxima) in the given time series by a straight line yields a continuous, piecewise linear function, $E_1(t)$. Similarily, connecting any two consecutive troughs (or local minima) in the time series by a straight line also yields a continuous, piecewise linear function, $E_2(t)$. Let the arithmetic mean of the above two functions be denoted by $E(t)$; i.e.,

$$E(t) = \frac{1}{2} (E_1(t) + E_2(t)) \tag{36}$$

and the oscillation of the original time series about this mean be denoted by $y(t)$; i.e.,

$$y(t) = x(t) - E(t). \tag{37}$$

It is obvious that $y(t)$ has a more rapid variation than $E(t)$ (Figure 1).  Thus, for the sake of convenient discussion, $y(t)$ and $E(t)$ will be called, respectively, the rapid-varying and slow-varying components of the original data.

## 2.2    SOME CHARACTERISTICS OF THE ENVELOPE ALGORITHM

Since the formulation of the envelope algorithm is not in a form amendable for a general mathematical analysis, some characteristics of the algorithm will be discussed by examining how it will split several types of simple deterministic stationary time series as given below.

(A)  For a time series consisting of two sine waves with different frequencies,

$$x(t) = a_1 \sin f_1 t + a_2 \sin f_2 t$$

with $f_1 < f_2$ and $a_1 f_1 < a_2 f_2$.  The slow-varying component will be, essentially,

$$E(t) = a_1 \sin f_1 t$$

and the rapid-varying component

$$y(t) = a_2 \sin f_2 t$$

That is to say, the envelope algorithm, essentially, breaks the sum of two sine waves into its individual wave form.  If the condition $a_1 f_1 < a_2 f_2$ is not met, then more complicated splitting than the above may result.

(B)  For a time series like an amplitude modulated signal, i.e.,

$$x(t) = (a + b \cos f_2 t) \cos f_1 t \qquad\qquad \text{with } f_1 > f_2 ,$$

where $f_1$ is the carrier frequency.  The slow-varying component will be, essentially, zero

$$E(t) = 0,$$

while the rapid-varying component is the same as the original signal. This is because both $E_1(t)$ and $E_2(t)$ are symmetric about $x = 0$ and with opposite signs.

(C)  For a time series like a frequency modulated signal; i.e.,

$$x(t) = a \cos[f_1 + f_2(t)]t,$$

where $f_1$ is the carrier frequency and $f_2(t)$ is the time dependent true signal. The slow-varying component is identical to the original signal. This is because both $E_1(t)$ and $E_2(t)$ are parallel to $x = 0$ with equal (but opposite) distance from $x = 0$.

(D)  Consider a more complex stationary deterministic time series defined as follows:

$$x(t) = a_1 \sin f_1 t + a_2 \sin f_2 t, \text{ for } 2mT \leq t < (2m + 1)T$$

$$= a_1 \sin f_1 t + a_3 \sin f_3 t, \text{ for } (2m + 1)T \leq t < (2m + 2)T$$

where $m = 0, 1, 2, 3$ and $T$ is some constant, with $f_1 < f_2 < f_3$, $a_1 f_1 < a_2 f_2$ and $a_1 f_1 < a_3 f_3$.

The corresponding slow-varying component will be, essentially,

$$E(t) = a_1 \sin f_1 t \quad \text{for } t \geq 0,$$

while the rapid-varying component will be

$$y(t) = a_2 \sin f_2 t \quad \text{for } 2mT \leq t < (2m+1)T$$

$$= a_3 \sin f_3 t \quad \text{for } (2m+1)T \leq t < (2m+2)T$$

This example points out that the envelope algorithm can cause different types of signal splitting on different portions of a time series, if signal wave form characteristics are different on different portions of the time series.

So far, only some deterministic stationary time series have been discussed. Of course, what is really of interest is the random time series, either stationary or non-stationary.  For a stationary random time series

$$x(t) = x_1(t) + x_2(t),$$

where $x_1(t)$ and $x_2(t)$ are both band-width limited white noises, respectively, with the following frequency ranges:

$$f_1 \leq f \leq f_2 \qquad \text{for } x_1(t)$$

$$f_3 \leq f \leq f_4 \qquad \text{for } x_2(t)$$

and

$$f_2 << f_3$$

The slow-varying and rapid-varying components will be, essentially,

$$E(t) = x_1(t)$$

and

$$y(t) = x_2(t),$$

for case (A).  On the other hand, for a stationary random time with arbitrary spectral density, there is simply no way to tell in advance what the slow varying and rapid-varying components will be.

## 2.3  APPLICATION TO THE SUNSPOT DATA

The record of sunspot number variation over the last 220 years which is often simply called the sunspot data, has been investigated extensively (refs. 9 through 11).  These data were chosen to test the performance of the envelope algorithm for detecting periodicity in the variation of the sunspot number.

The sunspot data used is a record of monthly average number of dark spots
on the sun surface from 1745 to 1965 (Figure 2).  Basically, it has been observed
that there is a dominant period of 22 years, which consists of one positive and
one negative phase with 11 years each, and that there might exist longer
periods - about 80 and 100 years.  However, these longer periods are not well
established yet, because the total available observation record is not suffi-
ciently long enough for an accurate statistical estimation.

The purpose of the test is to find out whether the above observed results
can be substantiated with the aid of the envelope algorithm.

Figures 3 and 4 show the slow-varying and rapid-varying components,
respectively, of the original data through one application of the envelope
algorithm.  This pair of components will be called as the first generation.
One can see clearly the justification for calling one component as the slow-
varying and the other as a rapid-varying component of the original time series.
Next, the obtained slow-varying component (Figure 3) is regarded as the original
time series.  The envelope algorithm is applied on it once again.  The second-
generation slow-varying and rapid-varying components are obtained as shown in
Figures 5 and 6.  The above process is repeated obtaining the third, fourth,
fifth, and higher-order generation components.  This process has been carried
through the fifth generation on the sunspot data.  They are shown in Figures 7
through 12.  The application of the envelope algorithm was stopped at the fifth
generation for the obvious reason that any further splitting of the resulting
slow-varying component will be good only for studying periodicity well exceeding
100 years.

Comparing both the slow-varying and rapid-varying components, respectively,
from the successive generations, the higher frequency components were gradually
removed by each application of the envelope algorithm as expected.  This is
analogous to successive applications of a series of low-pass filters, each of
which has progressively higher low cut-off frequency, on the resulting output
from the preceding filtering operation.  The advantage of the envelope algorithm
technique over the conventional digital filtering technique is its faster
computation and no requirement for specifying any integration time constant,

while the disadvantage is the lack of information on the frequency characteristics for each stage of signal splitting.

In order to investigate the periodicity of the sunspot data, the autocorrelation coefficient for the original data and the rapid-varying component from each generation were calculated. Figures 13 and 14 show the autocorrelations for the original data with integration constants of 110 years and 78 years, respectively. Both curves indicates very clearly the periodicity of about 11 years. In addition, Figure 14 may also indicate a period of 80 years (960 months).

Figures 15 through 20 show the autocorrelation coefficients for the rapid-varying components from the $1^{st}$ through $5^{th}$ generations by the envelope algorithm. The first five autocorrelations used an integration constant of 110 years, while the last curve used 78 years. No obvious peak except at t=0 is seen in the first-generation autocorrelation. In the second generation autocorrelation, a possible period of 11 years started to emerge, and a more definite indication of this period appears in the third generation autocorrelation curve. No clear and definite indication of periods may be derived from the fourth generation autocorrelation. The fifth generation autocorrelations with T=110 and 78 years both indicate a very definite period of about 105 years.

It may be concluded that the periodicity of the sunspot data of 11 years, 80 years, and 100 years, which are reported by previous investigations, has been confirmed again by the correlation analysis of the components obtained by the envelope algorithm. It should be noted that these different periods were observed from the different generations, but there is no rule established for selecting the number of generations one should apply to the envelope algorithm. Further tests of the envelope algorithm on various types of time series will be needed before more definite performance of the new algorithm for splitting time series may be known.

# Section III

# CONCLUSIONS

It has been shown that the crossed-beam test data obtained in the atmosphere are highly nonstationary and that for some cases the piecewise detrending and weighting methods can be employed to suppress the nonstationary trends to satisfactory results. However, it was also found from more recent analysis of infrared crossed-beam data that the above piecewise techniques were not completely adequate to deal with the nonstationarity encountered in strong environmental variations.

From the review of the available existing methods for analysis of non-stationary random processes, it is clear that none of the existing methods is generally applicable to all nonstationary random time series. As for the crossed-beam data, it seems that a logical way is to try the trend-elimination approach as described in Section I. It is thus recommended to implement a computer program for the algorithms, namely, equations (22) through (24) incorporating Tables 1 and 2. Secondly, the program will be used to process the Haswell crossed-beam test data (ref. 12). Finally, the resulting crossed-beam prediction will then be compared with other independent direct measurements for establishing the feasibility of suppressing nonstationary environmental variations by the trend elimination method.

Some characteristics of the envelope algorithm for splitting the random time series have been discussed based on several types of stationary time series. The actual application of the algorithm to the sunspot data has confirmed the previously observed periods of 11, 80, and 100 years. The definite performance of the envelope algorithm is, however, still not certain; more tests on various types of random time series are warranted.

# Section IV

# REFERENCES

1. Bendat, J. S. and Piersol, A. G., Measurement and Analysis of Random Data, John Wiley & Sons, Inc., 1966.

2. Kendall, M. G. and Stuart, A., The Advanced Theory of Statistics, Griffith, Vol. III, 1966.

3. Wiereville, W. W., "A Theory and Method For Correlation Analysis of Non-stationary Signals", IEEE Trans. Vol. EC-14, No. 6, December 1965.

4. Jayroe, R. R. and Su, M. Y., "Optimum Averaging Times of Meteorological Data With Time Dependent Means", NASA TMX-53782, October 1968.

5. Krause, F. R. and Hablntzel, B. C., "Noise Elimination by Piecewise Cross-Correlation of Photometer Outputs", Proc. Scientific Meeting of the Panel on Remote Atmospheric Probing, National Academy of Sciences, 16-17 May 1968.

6. Krause, F. R., Jones, J. A., Fisher, M. J., and Pooley, J. C., "Digital Analysis of Random Data Records by Piecewise Accumulation of Time Averages", NASA TND-6073, December 1970.

7. Pugachev, V. A., Theory of Random Function, Chapter 15, Pergaman Press, 1965.

8. Deutsch, R., Nonlinear Transformation of Random Processes, Prentice-Hall, 1962.

9. Fritz, H., "The Periods of Solar and Terrestrial Phenomena", Mon. Weath. Rev. Vol. 56. PP. 401-407, 1929.

10. Anderson, C. N., "Notes on the Sunspot Cycle", Jour. Geophysical Research. Vol. 59. No. 4, December 1954, pp. 455-461.

11. Sleeper, H. P., Jr., "Bi-stable Oscillation Modes of the Sun and Long Range Prediction of Solar Activity", Northrop Technical Report TR-793-709, February 1970.

12. Krause, F. R., Derr, V. E., et al., "Remote Probing of Wind and Turbulence Through Cross-Correlation of Passive Signals", Proc. 6th International Symp. Remote Sensing of Environment", 13-16 October 1969. PP. 327-359.

NOTE: ALL THREE DATA LINES ARE COMPOSED OF STRAIGHT SEGMENTS.

Figure 1.   AN ILLUSTRATION OF THE ENVELOPE ALGORITHM

Figure 2.  MONTHLY RECORD OF SUNSPOT NUMBER VARIATION



Figure 3.  SLOWING-VARYING COMPONENT OF SUNSPOT DATA BY ONE APPLICATION OF THE ENVELOPE ALGORITHM
(FIRST GENERATION)

Figure 4.  RAPID-VARYING COMPONENT OF SUNSPOT DATA BY ONE APPLICATION OF THE ENVELOPE ALGORITHM
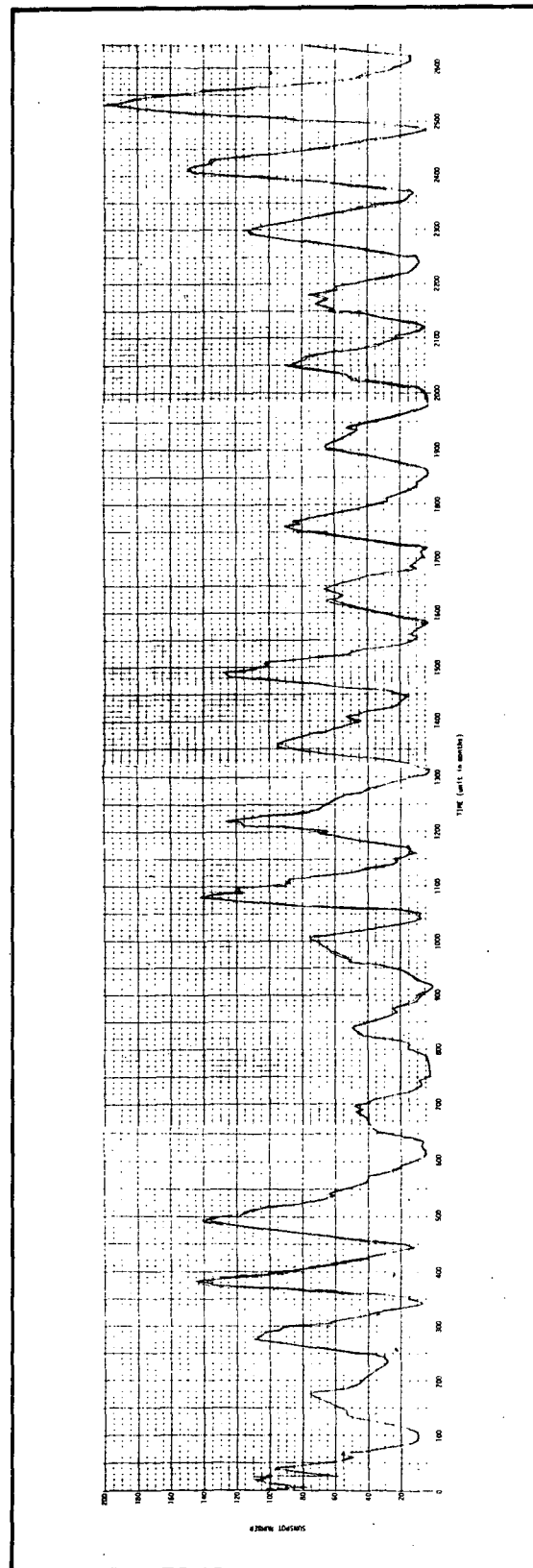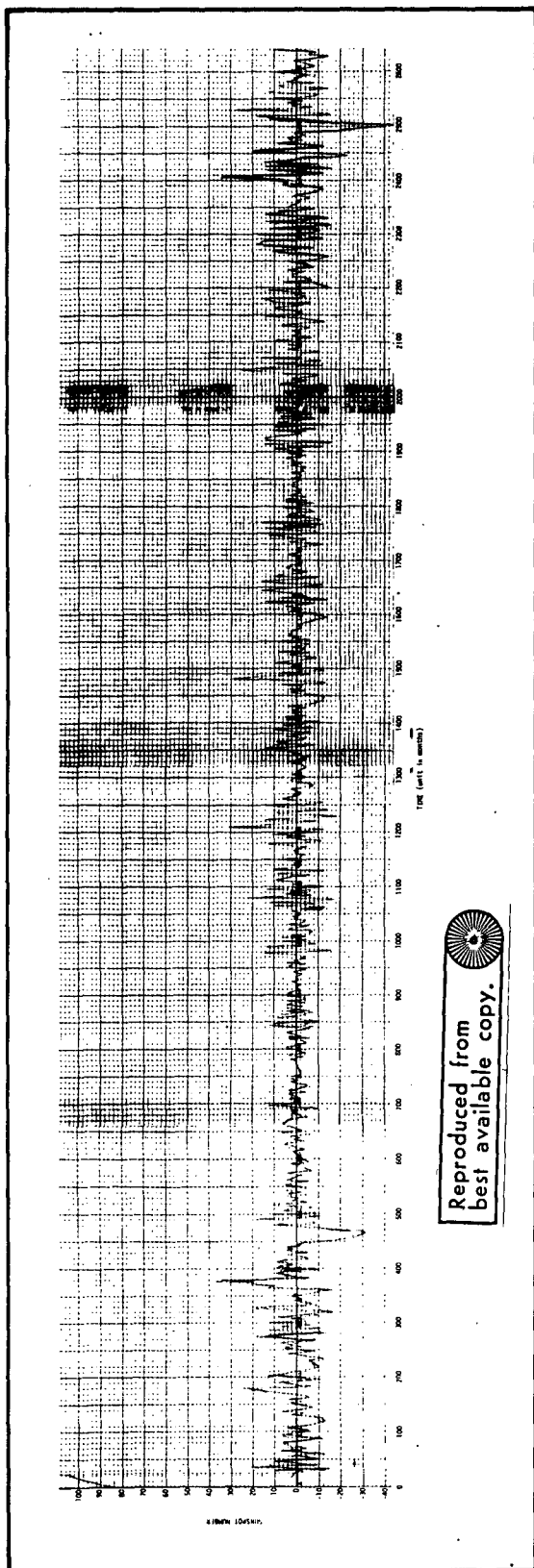(FIRST GENERATION)



Figure 5.  THE SECOND-GENERATION SLOW-VARYING COMPONENT OF SUNSPOT DATA

Figure 6.   THE SECOND-GENERATION RAPID-VARYING COMPONENT OF SUNSPOT DATA



Figure 7.   THE THIRD-GENERATION SLOW-VARYING COMPONENT OF SUNSPOT DATA

Figure 8.   THE THIRD-GENERATION RAPID-VARYING COMPONENT OF SUNSPOT DATA



Figure 9.   THE FOURTH-GENERATION SLOW-VARYING COMPONENT OF SUNSPOT DATA

Figure 10. THE FOURTH-GENERATION RAPID-VARYING COMPONENT OF SUNSPOT DATA



Figure 11. THE FIFTH-GENERATION SLOW-VARYING COMPONENT OF SUNSPOT DATA

Figure 12.  THE FIFTH-GENERATION RAPID-VARYING COMPONENT OF SUNSPOT DATA



Figure 13.  AUTOCORRELATION COEFFICIENT OF THE ORIGINAL SUNSPOT DATA WITH AN
INTEGRATION CONSTANT OF 110 YEARS

Figure 14. AUTOCORRELATION COEFFICIENT OF THE ORIGINAL SUNSPOT DATA WITH AN
INTEGRATION CONSTANT OF 110 YEARS



Figure 15. AUTOCORRELATION COEFFICIENT OF THE FIRST-GENERATION RAPID-VARYING
COMPONENT OF SUNSPOT DATA WITH INTEGRATION CONSTANT OF 110 YEARS

Figure 16. AUTOCORRELATION COEFFICIENT OF THE SECOND-GENERATION RAPID-VARYING
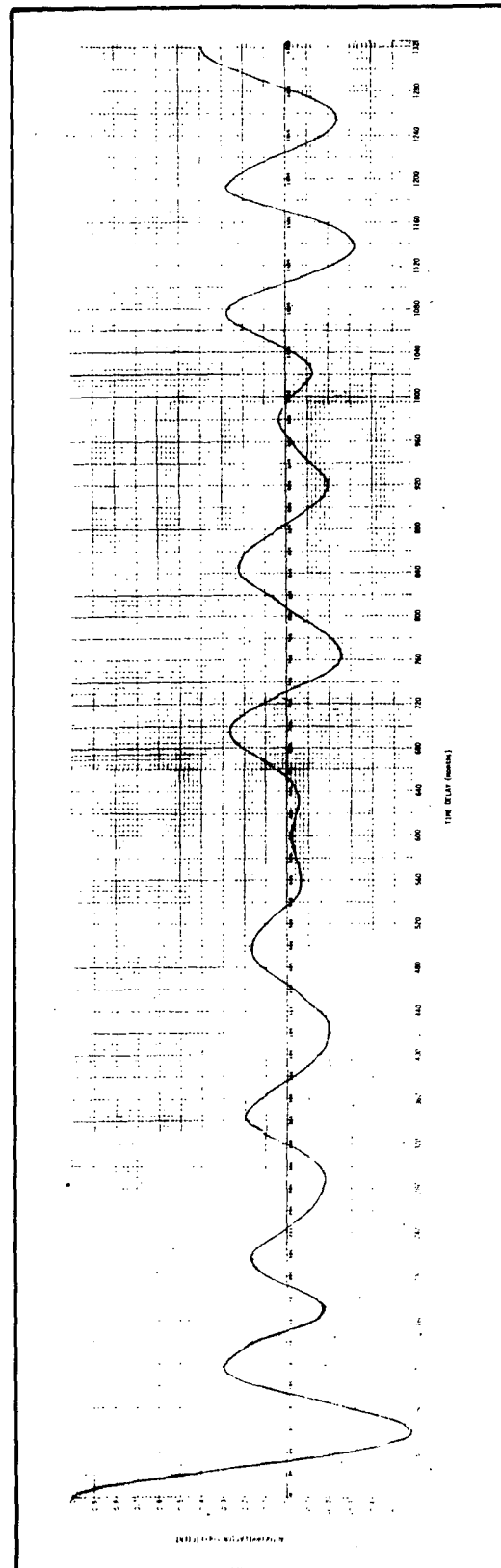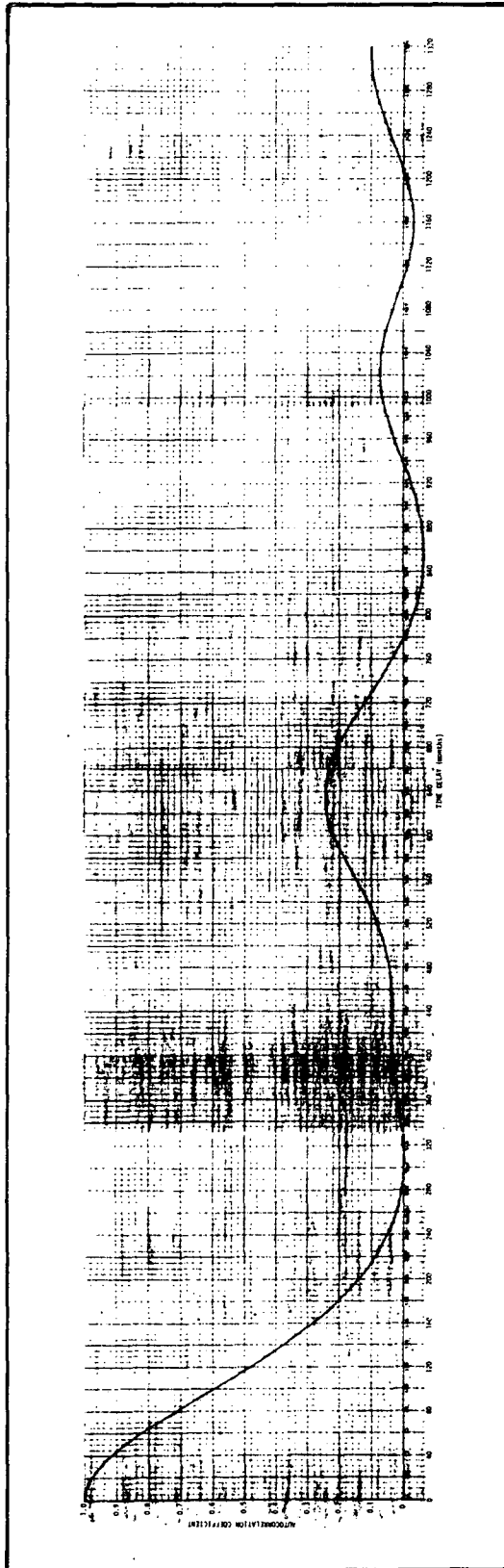COMPONENT OF SUNSPOT DATA WITH INTEGRATION CONSTANT OF 110 YEARS



Figure 17. AUTOCORRELATION COEFFICIENT OF THE THIRD-GENERATION RAPID-VARYING
COMPONENT OF SUNSPOT DATA WITH INTEGRATION CONSTANT OF 110 YEARS

Figure 18. AUTOCORRELATION COEFFICIENT OF THE FOURTH-GENERATION RAPID-VARYING COMPONENT OF SUNSPOT DATA WITH INTEGRATION CONSTANT OF 110 YEARS
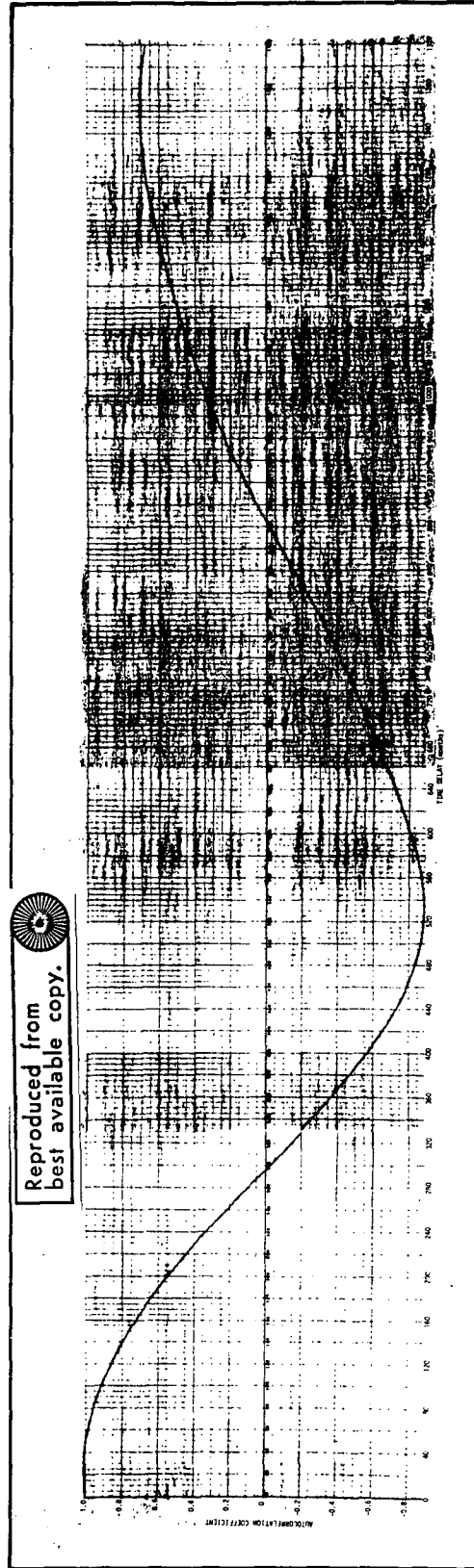
Figure 19. AUTOCORRELATION COEFFICIENT OF THE FIFTH-GENERATION RAPID-VARYING COMPONENT OF SUNSPOT DATA WITH INTEGRATION CONSTANT OF 110 YEARS
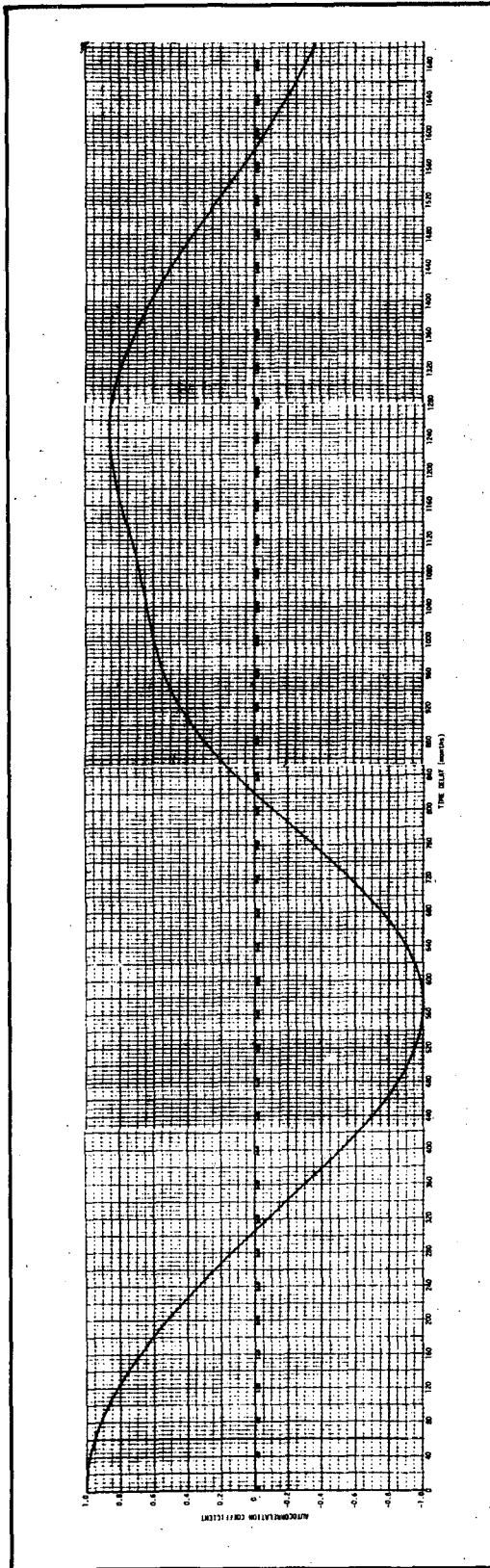
Figure 20.   AUTOCORRELATION COEFFICIENT OF THE FIFTH-GENERATION RAPID-VARYING COMPONENT
OF SUNSPOT DATA WITH INTEGRATION CONSTANT OF 78 YEARS